Introduction to Imitation Learning

Matt Barnes

TAs: Matthew Rockett, Gilwoo Lee, Matt Schmittle

*Some content borrowed from Florian Shkurti, Yisong Yue and Hoang M. Le

Recap

• Markov Decision Processes are a very general class of models, which encompass planning and reinforcement learning.



"Markov" means that _____ captures all information about the history x_1, x_2, \ldots, x_t

The most recent state x_t

Recap

_____ are known

transition model / dynamics / environment

Recap

• The difference between planning and reinforcement learning is whether the

Recap

• The three general methods for reinforcement learning are... (1) Model-based (2) Approximate dynamic programming

• (2) and (3) are both _____ methods

- (3) Policy gradient

 - Model-free

The multi-armed bandit problem is reinforcement learning with $__$ state(s)

1

Recap

problems?

Exploration vs. exploitation

Recap

What is the fundamental trade-off in bandit (and reinforcement learning)

Reward is equivalent to _____

Negative cost

Recap

The ε -greedy algorithm randomly explores with probability

 ${\cal E}$

Recap

bonus (i.e. confidence interval) which decreases with respect to $_____$

The number of times we try that action.

Recap

The UCB algorithm chooses actions according to the estimated reward, plus a

Today's lecture



Today's lecture

- from scratch.
- them.
- their own.

• Robots do not operate in a vacuum. They do not need to learn everything

• Humans need to easily interact with robots and share our expertise with

• Robots need to learn from the behavior and experience of others, not just

*Based off Florian Shkurti's lectures



Today's lecture

• Part 1: How can robots easily understand our objectives from demonstrations?

• Part 2: How can robots incorporate other's decision into their own?

Inverse reinforcement learning (IRL)

Behavior cloning (BC), interactive imitation learning (IL)

*Based off Florian Shkurti's lectures



Outp

		Policy Learning
Part 1	Inverse RL	No
	Behavior Cloning	Yes (direct)
Part 2	Interactive IL	Yes (direct)
	GAIL	Yes (indirect)

Today's lecture

out		Inputs	
Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstration
Yes	Yes	No	Yes
No	No	No	Yes
No	Yes	Yes	(Optional)
No	Yes	No	Yes

*Partially based off slides from Yisong Yue and Hoang M. Le



Part 1: Inverse Reinforcement Learning (IRL)

- Setting: No reward function. For complex tasks, these can be hard to specify!
- Fortunately, we are given a set of demonstrations

$$D = \{\tau_1, \dots, \tau_m\} = \prod_{\substack{m \text{ trajectories}}} T_{m}$$

• Goal: Learn a reward function r^* such that

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[r^*(x, u) \right]$$

$$\{(x_0^i, u_0^i, x_1^i, u_1^i...)\} \sim \rho_{\pi^*}$$

The state-action distribution of policy π^*

The high-level recipe

- Step 1: Learn a policy for the current reward function r
- Step 2: Update the reward function
- Repeat until policy trajectories are similar to demonstrations



Learning a reward function is an under defined problem

• Many reward functions correspond to the same policy



*Partially based off slides from Yisong Yue and Hoang M. Le



Learning a reward function is an under defined problem

- Let r^* be one solution, i.e. $\pi^* = \mathbf{a}$
- Note that ar^* for any constant a is also a solution, since $\operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ar^{*}(x, u) \right]$
- In fact, $r^* = 0$ is always a solution, since every policy is optimal

$$\operatorname{nrgmax}_{\pi} \mathbb{E}_{\pi} \left[r^*(x, u) \right]$$

$$= \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[r^*(x, u) \right]$$

Many different IRL approaches

- The reward function is linear [Abbeel & Ng 2004]
- [Ziebart et al., 2008]
- Maximum Margin Planning [Ratliff et al., 2006]

• Maximize the trajectory entropy, subject to a feature matching constraint

The linear reward function assumption reduces to feature expectation matching

• Assume $r(x) = \theta \cdot \phi(x)$ where $\phi(x)$ are the features of state x

• Then the value of a policy is linear in the expected features



The expected features of the policy

New objective, matching the feature expectations

- feature expectation μ
- Step 2: Update the reward function

Repeat until feature expectations are close $||\mu_{\pi} - \mu_{\pi^*}||$



• Step 1: Learn a policy π for the current reward function and compute its

$$\inf_{\substack{||\theta|| \le 1}} \max_{\theta} \theta^T (\mu_{\pi} - \mu_{\pi^*})$$

Update reward **Run RL** 21

*Partially based off slides from Yisong Yue and Hoang M. Le



Using IRL to predict pedestrian intention





[Kitani et al., 2012]

Part 2: Directly learning a policy

		Output		Inputs		
		Policy Learning	Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstration
Part 1	Inverse RL	No	Yes	Yes	No	Yes
	Behavior Cloning	Yes (direct)	No	No	No	Yes
Part 2	Interactive IL	Yes (direct)	No	Yes	Yes	(Optional)
	GAIL	Yes (indirect)	No	Yes	No	Yes



Behavior cloning

• Observe **pre-collected** expert The demonstrations:

 $D = \{\tau_1, \dots, \tau_m\} = \{(x_0^i, u_0^i, x_1^i, u_1^i \dots)\} \sim \rho_{\pi^*}$

• Learn a policy



The state-action distribution of policy π^{\ast}



Behavior cloning suffers from compounding errors

- The good news: $\hat{\pi}$ will perform well on samples from ρ_{π^*}
- The really bad news: When we roll out policy $\hat{\pi}$ it will inevitably make some mistakes compared to π^* , and these errors could compound resulting in drastically different state action distributions ρ_{π^*} and $\rho_{\hat{\pi}}$



Behavior cloning suffers from compounding errors

- The good news: $\hat{\pi}$ will perform well on samples from ρ_{π^*}
- The really bad news: When we roll out policy $\hat{\pi}$ it will inevitably make some mistakes compared to π^* , and these errors could compound resulting in drastically different state action distributions ρ_{π^*} and $\rho_{\hat{\pi}}$



Behavior cloning suffers from compounding errors

Theorem (simplified) [Ross et al., 2011]. Let ε be the supervised learning error rate of $\hat{\pi}$. Then the cumulative reward of this policy is bounded by:

Errors compound with respect to the time horizon T

 $J(\hat{\pi}) \le J(\pi^*) + T^2 \epsilon$

A history of covariate shift in imitation learning



30 x 32 pixels, 3 layer network, outputs steering command from approximately 5 minutes of training data per road type [Pomerleau 1992]



Navlab 1 (1986-1989) and Navlab 2 + ALVINN

A history of covariate shift in imitation learning



Figure 3.4: The single original video image is shifted and rotated to create multiple raining exemplars in which the vehicle appears to be at different locations relative to he road.



[Pomerleau 1992]



Imitation Learning is not supervised learning

- Policy's actions affect future observations/data
- This is not the case in supervised learning

Imitation Learning

- Train/test data are not i.i.d.
- If expected hold-out error is ε , then expected test error after T decisions is up to

 $T^2\varepsilon$

Errors compound

Supervised Learning

- Train/test data are i.i.d.
- If expected hold-out error is ε , then expected test error after T decisions is order

 $T\varepsilon$

Errors are independent

*From Florian Shkurti



Interactive imitation learning

		Output		Inputs		
		Policy Learning	Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstration
Part 1	Inverse RL	No	Yes	Yes	No	Yes
	Behavior Cloning	Yes (direct)	No	No	No	Yes
Part 2	Interactive IL	Yes (direct)	No	Yes	Yes	(Optional)
	GAIL	Yes (indirect)	No	Yes	No	Yes

*Partially based off slides from Yisong Yue and Hoang M. Le



Interactive feedback

- Roll-out any policy, and expert provides feedback for the current state
- In today's lecture, we'll consider the simple setting where this feedback takes the form of the expert's action π^{*}(x_t)

Expert feedback $\pi^*(x_t)$

The high-level recipe

- visits
- Step 2: Update the dataset and retrain the policy

Repeat

• Step 1: Roll-out the current policy, collect expert feedback on the states it

Initialize T policies, π_1, \ldots, π_T

• Step 1: Roll-out policy $u_1 \sim \pi_1(s_1), u_2 \sim \pi_2(s_2), ...$

and collect expert feedback $u_1^* = \pi^*(x_1), u_2^* = \pi^*(x_2), \dots$

• Step 2: Update policies

Intuitively, each policy will have to learn to correct for the mistakes of earlier policies

The Forward Training algorithm

The Forward Training algorithm

the cumulative reward of this policy is bounded by:



The downside: We have to learn T separate policies

• Theorem (simplified) [Ross et al., 2011]. Let ε be the supervised learning error rate of $\hat{\pi}$. Then

 $J(\hat{\pi}) \le J(\pi^*) + uT\epsilon$

Errors increase linearly with respect to the time horizon T, Same as supervised learning!



DAgger: Dataset Aggregation

• Learn only a single policy

Algorithm 1 DAgger

- 1: $D = \{(s, a)\}$ initial expert demonstrations
- 2: $\theta_1 \leftarrow$ train learner's policy parameters on D
- 3: for i = 1...N do
- 4:
- 5:
- 6:
- Train learner's policy $\pi_{\theta_{i+1}}$ on dataset D 7:
- 8: Return one of the policies π_{θ_i} that performs best on validation set

Execute learner's policy π_{θ_i} , get visited states $S_{\theta_i} = \{s_0, ..., s_T\}$ Query the expert at those states to get actions $A = \{a_0, ..., a_T\}$ Aggregate dataset $D = D \cup \{(s, a) \mid s \in S_{\theta_i}, a \in A\}$ Aggregate the data

DAgger: Dataset Aggregation





Ross et al, 2011

DAgger: Dataset Aggregation



Ross et al, 2011

GAIL: Generative Adversarial Imitation Learning

		Output		Inputs		
		Policy Learning	Reward Learning	Access to Environment	Interactive Demonstrator	Pre-collected Demonstration
Part 1	Inverse RL	No	Yes	Yes	No	Yes
	Behavior Cloning	Yes (direct)	No	No	No	Yes
Part 2	Interactive IL	Yes (direct)	No	Yes	Yes	(Optional)
	GAIL	Yes (indirect)	No	Yes	No	Yes

GAIL: Generative Adversarial Imitation Learning

• Setting: Pre-collected expert demonstrations

$$D = \{\tau_1, \dots, \tau_m\} = \{(x_0^i, u_0^i, x_1^i, u_1^i \dots)\} \sim \rho_{\pi^*}$$

• Goal: Minimize the *divergence* between ρ_{π^*} and $\rho_{\hat{\pi}}$

 $\hat{\pi} = \operatorname{argmin}_{\pi} D(\rho_{\pi} | | \rho_{\pi^*})$

P(x)

• Tells us how far apart two *distributions* are

[Nguyen et al., 2008])

What's a divergence?



Х

• Given samples from ρ_{π^*} and $\rho_{\hat{\pi}}$, we can estimate the divergence (e.g. using

GAIL: Generative Adversarial Imitation Learning

• The GAIL objective



- divergence
- Step 2: Update the policy using reinforcement learning.

• Step 1: For the current π , maximize the discriminator to estimate the



Takeaways

- speed up policy training.
- Forward Training or DAgger, which use interactive expert feedback.
- Inverse reinforcement learning allows us to learn an expert's reward function.

• Expert demonstrations, and in particular expert feedback, can dramatically

• Behavior Cloning suffers from compounding errors. This is not true for